

---

# Meta Flow Matching: Integrating Vector Fields on the Wasserstein Manifold

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Numerous biological and physical processes can be modeled as systems of interact-  
2 ing samples evolving continuously over time, e.g. the dynamics of communicating  
3 cells or physical particles. Flow-based models allow for learning these dynamics at  
4 the population level — they model the evolution of the entire distribution of sam-  
5 ples. However, current flow-based models are limited to a single initial population  
6 and a set of predefined conditions which describe different dynamics. We argue that  
7 multiple processes in natural sciences have to be represented as vector fields on the  
8 Wasserstein manifold of probability densities. That is, the change of the population  
9 at any moment in time depends on the population itself due to the interactions  
10 between samples. In particular, this is crucial for personalized medicine where the  
11 development of diseases and their treatments depend on the microenvironment of  
12 cells specific to each patient. We propose *Meta Flow Matching* (MFM), a practical  
13 approach to integrating along these vector fields on the Wasserstein manifold by  
14 amortizing the flow model over the initial populations. Namely, we embed the  
15 population of samples using a Graph Neural Network (GNN) and use these embed-  
16 dings to train a *Flow Matching* model. This gives Meta Flow Matching the ability  
17 to generalize over the initial distributions unlike previously proposed methods.  
18 Finally, we demonstrate the ability of MFM to improve prediction of individual  
19 treatment responses on a large scale multi-patient single-cell drug screen dataset.

## 20 1 Introduction

21 Understanding the dynamics of many-body problems is a central challenge across the natural sciences.  
22 In the field of cell biology, a central focus is the understanding of the dynamic processes that cells  
23 undergo in response to their environment, and in particular their response and interaction with other  
24 cells. Cells communicate with one other in close proximity using *cell signaling*, exerting influence  
25 over each other’s trajectories (Armingol et al., 2020; Goodenough and Paul, 2009). This signaling  
26 presents an obstacle for modeling, but is essential for understanding and eventually controlling  
27 cell dynamics during development (Gulati et al., 2020; Rizvi et al., 2017), in diseased states (Molè  
28 et al., 2021; Binnewies et al., 2018; Zeng and Dai, 2019; Chung et al., 2017), and in response to  
29 perturbations (Ji et al., 2021; Peidli et al., 2024).

30 The super-exponential decrease of sequencing costs and advances in microfluidics has enabled the  
31 rapid advancement of single-cell sequencing and related technologies over the past decade. While  
32 single-cell sequencing has been used to great effect to understand the heterogeneity in cell systems,  
33 they are also destructive, making longitudinal measurements extremely difficult. Instead, most  
34 approaches model cell dynamics at the population level (Hashimoto et al., 2016; Weinreb et al., 2018;  
35 Schiebinger et al., 2019; Tong et al., 2020; Neklyudov et al., 2022; Bunne et al., 2023a). These  
36 approaches involve the formalisms of optimal transport (Villani, 2009; Peyré and Cuturi, 2019) and

37 generative modeling (De Bortoli et al., 2021; Lipman et al., 2023) methods, which allow for learning  
 38 a map between empirical measures. While these methods are able to model the dynamics of the  
 39 population, they are fundamentally limited in that they model the evolution of cells as independent  
 40 particles evolving according to a shared dynamical system. Furthermore, these models can be trained  
 41 to match any given set of measures, but they are restricted to modeling of a single population and can  
 42 at best condition on a number of different dynamics that is available in the training data.

43 To address this we propose *Meta Flow Matching* (MFM) — the amortization of the Flow Matching  
 44 generative modeling framework (Lipman et al., 2023) over the input measures. In practice, our  
 45 method can be used to predict the time-evolution of distributions from a given dataset of the time-  
 46 evolved examples. Namely, we assume that the collected data undergoes a universal developmental  
 47 process, which depends only on the population itself as in the setting of the interacting particles or  
 48 communicating cells. Under this assumption, we learn the vector field model that takes samples from  
 49 the initial distribution as input and defines the push-forward map on the sample-space that maps the  
 50 initial distribution to the final distribution.

51 We showcase the utility of our approach on two applications. We first explore Meta Flow Matching on  
 52 a synthetic task of denoising letters. We show that MFM is able to generalize the denoising process  
 53 to letters in unseen orientations where a standard flow matching approach cannot. Next, we explore  
 54 how MFM can be applied to model single-cell perturbation data (Ji et al., 2021; Peidli et al., 2024).  
 55 We evaluate MFM on predicting the response of patient-derived cells to chemotherapy treatments  
 56 in a recently published large scale single-cell drug screening dataset where there are known to be  
 57 patient-specific responses (Ramos Zapatero et al., 2023). This dataset includes more than 25M cells  
 58 collected over ten patients under 2500 conditions. This is a challenging task due to the variance over  
 59 multiple patients, treatments applied and the local cell compositions, but it can be used to study the  
 60 *tumor micro-environment* (TME), thought to be essential in circumventing chemoresistance. We  
 61 demonstrate that Meta Flow Matching can successfully predict the development of cell populations  
 62 on replicated experiments, and, most importantly, it generalizes to previously unseen patients, thus,  
 63 capturing the patient-specific response to the treatment.

## 64 2 Background

### 65 2.1 Generative Modeling via Flow Matching

66 Flow Matching is an approach to generative modeling recently proposed independently in different  
 67 works: Rectified Flows (Liu et al., 2022), Flow Matching (Lipman et al., 2023), Stochastic Interpolants  
 68 (Albergo and Vanden-Eijnden, 2022). It assumes a continuous interpolation between densities  $p_0(x_0)$   
 69 and  $p_1(x_1)$  in the sample space. That is, the sample from the intermediate density  $p_t(x_t)$  is produced  
 70 as follows

$$x_t = f_t(x_0, x_1), \quad (x_0, x_1) \sim \pi(x_0, x_1), \quad (1)$$

$$\text{where } \int dx_1 \pi(x_0, x_1) = p_0(x_0), \quad \int dx_0 \pi(x_0, x_1) = p_1(x_1), \quad (2)$$

71 where  $f_t$  is the time-continuous interpolating function such that  $f_{t=0}(x_0, x_1) = x_0$  and  
 72  $f_{t=1}(x_0, x_1) = x_1$  (e.g. linearly between  $x_0$  and  $x_1$  with  $f_t(x_0, x_1) = (1-t) \cdot x_0 + t \cdot x_1$ );  
 73  $\pi(x_0, x_1)$  is the density of the joint distribution, which is usually taken as a distribution of inde-  
 74 pendent random variables  $\pi(x_0, x_1) = p_0(x_0)p_1(x_1)$ , but can also be generalized to formulate the  
 75 optimal transport problems (Pooladian et al., 2023; Tong et al., 2024). The corresponding density can  
 76 be defined then as the following expectation

$$p_t(x) = \int dx_0 dx_1 \pi(x_0, x_1) \delta(x - f_t(x_0, x_1)). \quad (3)$$

77 The essential part of Flow Matching is the continuity equation that describes the change of this  
 78 density through the vector field on the state space, which admits vector field  $v_t^*(x)$  as a solution

$$\frac{\partial p_t(x)}{\partial t} = -\langle \nabla_x, p_t(x) v_t^*(x) \rangle, \quad v_t^*(\xi) = \frac{1}{p_t(\xi)} \mathbb{E}_{\pi(x_0, x_1)} \left[ \delta(f_t(x_0, x_1) - \xi) \frac{\partial f_t(x_0, x_1)}{\partial t} \right]. \quad (4)$$

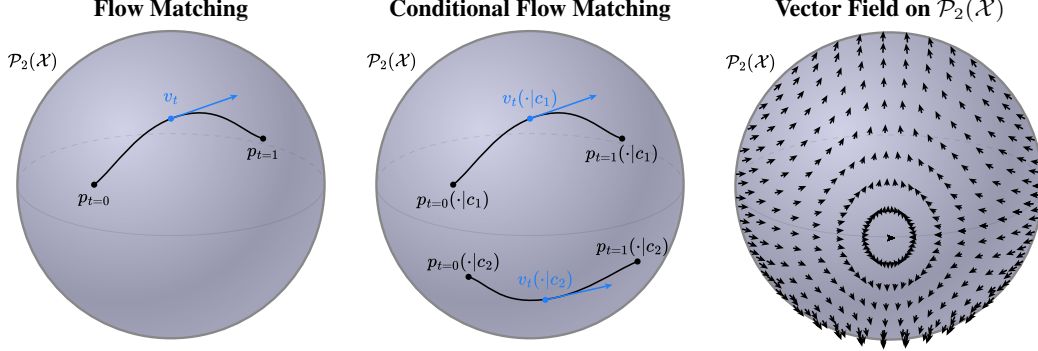


Figure 1: Illustration of flow matching methods on the 2-Wasserstein manifold,  $\mathcal{P}_2(\mathcal{X})$ , depicted as a two-dimensional sphere. *Flow Matching* learns the tangent vectors to a single curve on the manifold. *Conditional* generation corresponds to learning a finite set of curves on the manifold, e.g. classes  $c_1$  and  $c_2$  on the plot. *Meta Flow Matching* learns to integrate a vector field on  $\mathcal{P}_2(\mathcal{X})$ , i.e. for every starting density  $p_0$  Meta Flow Matching defines a push-forward measure that integrates along the underlying vector field.

79 Relying on this formula, one can derive the tractable objective for learning  $v_t^*(x)$ , i.e.

$$\mathcal{L}_{\text{FM}}(\omega) = \int_0^1 dt \mathbb{E}_{p_t(x)} \|v_t^*(x) - v_t(x; \omega)\|^2 \quad (5)$$

$$= \mathbb{E}_{\pi(x_0, x_1)} \int_0^1 dt \left\| \frac{\partial}{\partial t} f_t(x_0, x_1) - v_t(f_t(x_0, x_1); \omega) \right\|^2 + \text{constant}. \quad (6)$$

80 Finally, the vector field  $v_t(\xi, \omega) \approx v_t^*(\xi)$  defines the push-forward density that approximately matches  
 81  $p_{t=1}$ , i.e.  $T_{\#}p_0 \approx p_{t=1}$ , where  $T$  is the flow corresponding to vector field  $v_t(\cdot, \omega)$  with parameters  $\omega$ .

## 82 2.2 Conditional Generative Modeling via Flow Matching

83 Conditional image generation is one of the most common applications of generative models nowadays;  
 84 it includes conditioning on the text prompts (Saharia et al., 2022b; Rombach et al., 2022) as well  
 85 as conditioning on other images (Saharia et al., 2022a). To learn the conditional generative process  
 86 with diffusion models, one merely has to pass the conditional variable (sampled jointly with the data  
 87 point) as an additional input to the parametric model of the vector field. The same applies for the  
 88 Flow Matching framework.

89 Conditional Generative Modeling via Flow Matching is independently introduced in several works  
 90 (Zheng et al., 2023; Dao et al., 2023; Isobe et al., 2024) and it operates as follows. Consider a family  
 91 of time-continuous densities  $p_t(x_t | c)$ , which corresponds to the distribution of the following random  
 92 variable

$$x_t = f_t(x_0, x_1), \quad (x_0, x_1) \sim \pi(x_0, x_1 | c). \quad (7)$$

93 For every  $c$ , the density  $p_t(x_t | c)$  follows the continuity equation with the following vector field

$$v_t^*(\xi | c) = \frac{1}{p_t(\xi | c)} \mathbb{E}_{\pi(x_0, x_1)} \delta(f_t(x_0, x_1) - \xi) \frac{\partial f_t(x_0, x_1)}{\partial t}, \quad (8)$$

94 which depends on  $c$ . Thus, the training objective of the conditional model becomes

$$\mathcal{L}_{\text{CGFM}}(\omega) = \mathbb{E}_{p(c)} \mathbb{E}_{\pi(x_0, x_1 | c)} \int_0^1 dt \left\| \frac{\partial}{\partial t} f_t(x_0, x_1) - v_t(f_t(x_0, x_1) | c; \omega) \right\|^2, \quad (9)$$

95 where, compared to the original Flow Matching formulation, we first have to sample  $c$ , then produce  
 96 the samples from  $p_t(x_t | c)$  and pass  $c$  as input to the parametric model of the vector field.

## 97 3 Meta Flow Matching

98 In this paper, we propose the amortization of the Flow Matching framework over the marginal  
 99 distributions. Our model is based on the outstanding ability of the Flow Matching framework to

100 learn the push-forward map for any joint distribution  $\pi(x_0, x_1)$  given empirically. For the given joint  
 101  $\pi(x_0, x_1)$ , we denote the solution of the Flow Matching optimization problem as follows

$$v_t^*(\cdot, \pi) = \underset{v_t}{\operatorname{argmin}} \mathcal{L}_{GFM}(v_t(\cdot), \pi(x_0, x_1)). \quad (10)$$

102 Analogously to the amortized optimization (Chen et al., 2022; Amos et al., 2023), we aim to learn the  
 103 model that outputs the solution of Eq. (10) based on the input data sampled from  $\pi$ , i.e.

$$v_t(\cdot, \varphi(\pi)) = v_t^*(\cdot, \pi), \quad (11)$$

104 where  $\varphi(\pi)$  is the embedding model of  $\pi$  and the joint density  $\pi(\cdot | c)$  is generated using some  
 105 unknown measure of the conditional variables  $c \sim p(c)$ .

### 106 3.1 Modeling Process in Natural Sciences as Vector Fields on the Wasserstein Manifold

107 We argue that numerous biological and physical processes cannot be modeled via the vector field  
 108 propagating the population samples independently. Thus, we propose to model these processes as  
 109 families of conditional vector fields where we amortize the conditional variable by embedding the  
 110 population via a Graph Neural Network (GNN).

111 To provide the reader with the necessary intuition, we are going to use the geometric formalism  
 112 developed by Otto (2001). That is, time-dependent densities  $p_t(x_t)$  define absolutely-continuous  
 113 curves on the 2-Wasserstein space of distributions  $\mathcal{P}_2(\mathcal{X})$  (Ambrosio et al., 2008). The tangent space  
 114 of this manifold is defined by the gradient flows  $\mathcal{S}_t = \{\nabla_{s_t} | s_t : \mathcal{X} \rightarrow \mathbb{R}\}$  on the state space  $\mathcal{X}$ . In  
 115 the Flow Matching context, we are going to refer to the tangent vectors as vector fields since one  
 116 can always project the vector field onto the tangent space by parameterizing it as a gradient flow  
 117 (Neklyudov et al., 2022).

118 Under the geometric formalism of the 2-Wasserstein manifold, Flow Matching can be considered  
 119 as learning the tangent vectors  $v_t(\cdot)$  along the density curve  $p_t(x_t)$  defined by the sampling process  
 120 in Eq. (2) (see the left panel in Fig. 1). Furthermore, the conditional generation processes  $p_t(x_t | c)$   
 121 would be represented as a finite set of curves if  $c$  is discrete (e.g. class-conditional generation of  
 122 images) or as a family of curves if  $c$  is continuous (see the middle panel in Fig. 1).

123 Finally, one can define a vector field on the 2-Wasserstein manifold via the continuity equation with  
 124 the vector field  $v_t(x, p_t(x))$  on the state space  $\mathcal{X}$  that depends on the current density  $p_t(x)$  or its  
 125 derivatives. Below we give two examples of processes defined as vector fields on the 2-Wasserstein  
 126 manifold.

127 **Example 1** (Mean-field limit of interacting particles). *In the limit of the infinite number of interacting*  
 128 *particles one can describe their state with the density function  $p_t(x)$ . Consider the interaction*  
 129 *according to the first order dynamics with the velocity  $k(x, y) : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}^d$  of the particles at*  
 130 *point  $x$  that interact with the particles at point  $y$ . Then the change of the density is described by the*  
 131 *following continuity equation*

$$\frac{dx}{dt} = \mathbb{E}_{p_t(y)} k(x, y), \quad \frac{\partial p_t(x)}{\partial t} = -\langle \nabla_x, p_t(x) \mathbb{E}_{p_t(y)} k(x, y) \rangle. \quad (12)$$

132 **Example 2** (Diffusion). *Even when the physical particles evolve independently in nature, the*  
 133 *deterministic vector field model might be dependent on the current density of the population. For*  
 134 *instance, for the diffusion process, the change of the density is described by the Fokker-Planck*  
 135 *equation, which results in the density-dependent vector field when written as a continuity equation,*  
 136 *i.e.*

$$\frac{\partial p_t(x)}{\partial t} = \frac{1}{2} \Delta_x p_t(x) = -\left\langle \nabla_x, p_t(x) \left( -\frac{1}{2} \nabla_x \log p_t(x) \right) \right\rangle \implies \frac{dx}{dt} = -\frac{1}{2} \nabla_x \log p_t(x). \quad (13)$$

137 Motivated by the examples above, we argue that using the information about the current or the initial  
 138 density is crucial for the modeling of time-evolution of densities in natural processes, to capture this  
 139 type of dependency one can model the change of the density as the following Cauchy problem

$$\frac{\partial p_t(x)}{\partial t} = -\langle \nabla_x, p_t(x) v_t(x, p_t) \rangle, \quad p_{t=0}(x) = p_0(x), \quad (14)$$

140 where the state-space vector field  $v_t(x, p_t)$  depends on the density  $p_t$ .

141 The dependency might vary across models, e.g. in Example 1 the vector field can be modeled as an  
 142 application of a kernel to the density function, while in Example 2 the vector field depends only on  
 143 the local value of the density and its derivative.

144 **3.2 Integrating Vector Fields on the Wasserstein Manifold via Meta Flow Matching**

145 Consider the dataset of joint populations  $\mathcal{D} = \{(\pi(x_0, x_1 | i))\}_i$ , where, to simplify the notation,  
 146 we associate every  $i$ -th population with its density  $\pi(\cdot | i)$  and the conditioning variable here is the  
 147 index of this population in the dataset. We make the following assumptions regarding the ground  
 148 truth sampling process (i) we assume that the starting marginals  $p_0(x_0 | i) = \int dx_1 \pi(x_0, x_1 | i)$  are  
 149 sampled from some unknown distribution that can be parameterized with a large enough number of  
 150 parameters (ii) the endpoint marginals  $p_1(x_1 | i) = \int dx_0 \pi(x_0, x_1 | i)$  are obtained as push-forward  
 151 densities solving the Cauchy problem in Eq. (14), (iii) there exists unique solution to this Cauchy  
 152 problem.

153 One can learn a joint model of all the processes from the dataset  $\mathcal{D}$  using the conditional version of  
 154 the Flow Matching algorithm (see Section 2.2) where the population index  $i$  plays the role of the  
 155 conditional variable. However, obviously, such a model will not generalize beyond the considered  
 156 data  $\mathcal{D}$  and unseen indices  $i$ . We illustrate this empirically in Section 5.

157 To be able to generalize to previously unseen populations, we propose learning the density-dependent  
 158 vector field motivated by Eq. (14). That is, we propose to use an embedding function  $\varphi : \mathcal{P}_2(\mathcal{X}) \rightarrow$   
 159  $\mathbb{R}^m$  to embed the starting marginal density  $p_0$ , which we then input into the vector field model and  
 160 minimize the following objective over  $\omega$

$$\mathcal{L}_{\text{MFM}}(\omega; \varphi) = \mathbb{E}_{i \sim \mathcal{D}} \mathbb{E}_{\pi(x_0, x_1 | i)} \int_0^1 dt \left\| \frac{\partial}{\partial t} f_t(x_0, x_1) - v_t(f_t(x_0, x_1) | \varphi(p_0); \omega) \right\|^2. \quad (15)$$

161 Note that the initial density  $p_0$  is enough to predict the push-forward density  $p_1$  since the Cauchy  
 162 problem for Eq. (14) has a unique solution. The embedding function  $\varphi(p_0)$  can take different forms,  
 163 e.g. it can be the density value  $\varphi(p_0) = p_0(\cdot)$ , which is then used inside the vector field model to  
 164 evaluate at the current point (analogous to Example 2); a kernel density estimator (analogous to  
 165 Example 1); or a parametric model taking the samples from this density as an input.

166 **Proposition 1.** *Meta Flow Matching recovers the Conditional Generation via Flow Matching*  
 167 *when the conditional dependence of the marginals  $p_0(x_0 | c) = \int dx_1 \pi(x_0, x_1 | c)$  and  $p_1(x_1 | c) =$*   
 168  *$\int dx_0 \pi(x_0, x_1 | c)$  and the distribution  $p(c)$  are known, i.e. there exist  $\varphi : \mathcal{P}_2(\mathcal{X}) \rightarrow \mathbb{R}^m$  such that*  
 169  *$\mathcal{L}_{\text{MFM}}(\omega) = \mathcal{L}_{\text{CGFM}}(\omega)$ .*

170 *Proof.* Indeed, sampling from the dataset  $i \sim \mathcal{D}$  becomes sampling of the conditional variable  
 171  $c \sim p(c)$  and the embedding function becomes  $\varphi(p_0(\cdot | c)) = c$ .  $\square$

172 Furthermore, for the parametric family of the embedding models  $\varphi(p_t, \theta)$ , we show that the parameters  
 173  $\theta$  can be estimated by minimizing the objective in Eq. (15) in the joint optimization with the vector  
 174 field parameters  $\omega$ . We formalize this statement in the following theorem.

175 **Theorem 1.** *Consider a dataset of populations  $\mathcal{D} = \{(\pi(x_0, x_1 | i))\}_i$  generated from some unknown*  
 176 *conditional model  $\pi(x_0, x_1 | c)p(c)$ . Then the following objective*

$$\mathcal{L}(\omega, \theta) = \mathbb{E}_{p(c)} \int_0^1 dt \mathbb{E}_{p_t(x_t | c)} \|v_t^*(x_t | c) - v_t(x_t | \varphi(p_0, \theta), \omega)\|^2 \quad (16)$$

177 *is equivalent to the Meta Flow Matching objective*

$$\mathcal{L}_{\text{MFM}}(\omega, \theta) = \mathbb{E}_{i \sim \mathcal{D}} \mathbb{E}_{\pi(x_0, x_1 | i)} \int_0^1 dt \left\| \frac{\partial}{\partial t} f_t(x_0, x_1) - v_t(f_t(x_0, x_1) | \varphi(p_0, \theta); \omega) \right\|^2 \quad (17)$$

178 *up to an additive constant.*

179 *Proof.* We postpone the proof to Appendix A.  $\square$

180 **3.3 Learning Population Embeddings via Graph Neural Networks (GNNs)**

181 In many applications, the populations  $\mathcal{D} = \{(\pi(x_0, x_1 | i))\}_{i=1}^N$  are given as empirical distributions,  
 182 i.e. they are represented as samples from some unknown density  $\pi$

$$\{(x_0^j, x_1^j)\}_{j=1}^{N_i}, \quad (x_0^j, x_1^j) \sim \pi(x_0, x_1 | i), \quad (18)$$

183 where  $N_i$  is the size of the  $i$ -th population. For instance, for the diffusion process considered in  
 184 Example 2, the samples from  $\pi(x_0, x_1 | i)$  can be generated by generating some marginal  $p_1(x_1 | i)$   
 185 and then adding the Gaussian random variable to the samples  $x_1^j$ . We use this model in our synthetic  
 186 experiments in Section 5.1.

187 Since the only available information about the populations is samples, we propose learning the  
 188 embedding of populations via a parametric model  $\varphi(p_0, \theta)$ , i.e.

$$\varphi(p_0, \theta) = \varphi\left(\{x_0^j\}_{j=1}^{N_i}, \theta\right), \quad (x_0^j, x_1^j) \sim \pi(x_0, x_1 | i). \quad (19)$$

189 For this purpose, we employ GNNs, which recently have been successfully applied for simulation of  
 190 complicated many-body problems in physics (Sanchez-Gonzalez et al., 2020). To embed a population  
 191  $\{x_0^j\}_{j=1}^{N_i}$ , we create a  $k$ -nearest neighbour graph  $G_i$  based on the metric in the state-space  $\mathcal{X}$ , input it  
 192 into a GNN, which consists of several message-passing iterations (Gilmer et al., 2017) and the final  
 193 average-pooling across nodes to produce the embedding vector. Finally, we update the parameters of  
 194 the GNN jointly with the parameters of the vector field to minimize the loss function in Eq. (17).

## 195 4 Related Work

196 The meta-learning of probability measures was previously studied by Amos et al. (2022) where they  
 197 demonstrate that the prediction of the optimal transport paths can be efficiently amortized over the  
 198 input marginal measures. The main difference with our approach is that we are trying to learn the  
 199 push-forward map without embedding the second marginal.

200 **Generative modeling for single cells.** Single cell data has expanded to encompass multiple modalities  
 201 of data profiling cell state and activities (Frangieh et al., 2021; Bunne et al., 2023b). Single-cell  
 202 data presents multiple challenges in terms of noise, non-time resolved, and high dimension, and  
 203 generative models have been used to counter those problems. Autoencoder has been used to embed  
 204 and extrapolate data Out Of Distribution (OOD) with its latent state dimension (Lotfollahi et al., 2019;  
 205 Lopez et al., 2018; Hetzel et al., 2022). Orthogonal non-negative matrix factorization (oNMF) has  
 206 also been used for dimensionality reduction combined with mixture models for cell state prediction  
 207 (Chen et al., 2020). Other approaches have tried to use Flow Matching (FM) (Tong et al., 2023, 2024;  
 208 Neklyudov et al., 2023) or similar approaches such as the Monge gap (Uscidda and Cuturi, 2023) to  
 209 predict cell trajectories. Currently, the state of the art method uses the principle of Optimal Transport  
 210 (OT) to predict cell trajectories with Input Convex Neural Network (ICNN) (Makkuva et al., 2020;  
 211 Bunne et al., 2023b). What determines the significance of the method is its capability in generalizing  
 212 out of distribution to a new population of cells, which may be from different culture or individuals.  
 213 As of this time, our method is the only method that takes inter-cellular interactions into account.

214 **Generative modeling for physical processes.** The closest approach to ours is the prediction of the  
 215 many-body interactions in physics (Sanchez-Gonzalez et al., 2020) via GNNs. However, the problem  
 216 there is very different since these models use the information about the individual trajectories of  
 217 samples, which are not available for the single-cell prediction. Neklyudov et al. (2022) consider  
 218 learning the vector field for any continuous time-evolution of a probability measure, however, their  
 219 method is restricted to single curves and do not consider generalization to unseen data. Finally, the  
 220 weather/climate forecast models generating the next state conditioned on the previous one (Price  
 221 et al., 2023; Verma et al., 2024) are similar approaches to ours but operating on a much finer time  
 222 resolution.

## 223 5 Experiments

224 To show the effectiveness of MFM to generalize under previously unseen populations for the task  
 225 population prediction, we consider two experimental settings. (i) A synthetic experiment with well  
 226 defined coupled populations, and (ii) experiments on a publicly available single-cell dataset consisting  
 227 of populations from patient dependent treatment response trials. To quantify model performance,  
 228 we consider three distributional distances metrics: the 1-Wasserstein distance ( $\mathcal{W}_1$ ), 2-Wasserstein  
 229 ( $\mathcal{W}_2$ ) distance, and the radial basis kernel maximum-mean-discrepancy (MMD) distance (Gretton  
 230 et al., 2012). We parameterize all vector field models  $v_t(\cdot | \varphi(p_0); \omega)$  using a Multi-Layer Perceptron  
 231 (MLP). For MFM, we additionally parameterize  $\varphi(p_t; \theta, k)$  using a Graph Convolutional Network

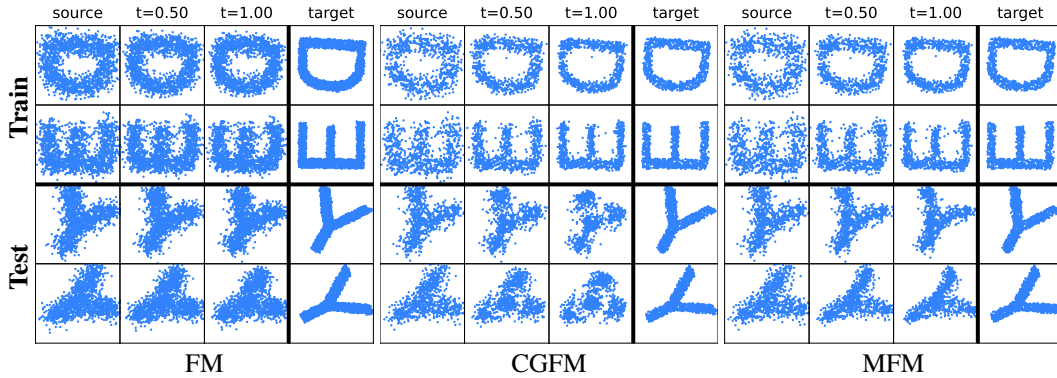


Figure 2: Examples of model-generated samples for synthetic letters from the source distribution ( $t = 0$ ) to predicted target distribution ( $t = 1$ ). See Fig. 4 in Appendix F for a larger set of examples.

Table 1: Results of the synthetic letters experiment for population prediction on seen train populations and unseen test populations. We report the the 1-Wasserstein ( $\mathcal{W}_1$ ), 2-Wasserstein ( $\mathcal{W}_2$ ), and the maximum-mean-discrepancy (MMD) distributional distances. We consider 4 settings for MFM with varying  $k$ .

	Train			Test		
	$\mathcal{W}_1$	$\mathcal{W}_2$	MMD ( $\times 10^{-3}$ )	$\mathcal{W}_1$	$\mathcal{W}_2$	MMD ( $\times 10^{-3}$ )
FM	$0.216 \pm 0.000$	$0.280 \pm 0.000$	$2.38 \pm 0.00$	$0.237 \pm 0.000$	$0.315 \pm 0.000$	<b><math>3.28 \pm 0.00</math></b>
CGFM	<b><math>0.093 \pm 0.000</math></b>	<b><math>0.112 \pm 0.000</math></b>	$0.34 \pm 0.00$	$0.317 \pm 0.000$	$0.397 \pm 0.000$	$6.67 \pm 0.00$
MFM ( $k = 0$ )	$0.099 \pm 0.000$	$0.128 \pm 0.000$	$0.25 \pm 0.00$	$0.221 \pm 0.000$	$0.267 \pm 0.000$	$3.77 \pm 0.00$
MFM ( $k = 1$ )	<b><math>0.096 \pm 0.003</math></b>	$0.124 \pm 0.004$	<b><math>0.22 \pm 0.04</math></b>	$0.217 \pm 0.003$	$0.261 \pm 0.003$	$3.80 \pm 0.28$
MFM ( $k = 10$ )	<b><math>0.096 \pm 0.003</math></b>	$0.124 \pm 0.003$	<b><math>0.23 \pm 0.04</math></b>	<b><math>0.213 \pm 0.008</math></b>	<b><math>0.256 \pm 0.008</math></b>	<b><math>3.68 \pm 0.45</math></b>
MFM ( $k = 50$ )	$0.099 \pm 0.003$	$0.127 \pm 0.003$	$0.25 \pm 0.05$	$0.226 \pm 0.005$	$0.270 \pm 0.007$	$4.38 \pm 0.30$

232 (GCN) with a  $k$ -nearest neighbour graph edge pooling layer. We include details regarding model  
 233 hyperparameters, training/optimization, and implementation in Appendix B and Appendix B.2. The  
 234 results for all the models are averaged over three random seeds.

## 235 5.1 Synthetic Experiment

236 We curate a synthetic dataset of the joint distributions  $\{(p_0(x_0, |i), p_1(x_1 | i))\}_{i=1}^N$  by simulating a  
 237 diffusion process applied to a set of pre-defined target distributions  $p_1(x_1 | i)$  for  $i = 1, \dots, N$ . To get  
 238 a paired population  $p_0(x_0 | i)$  we simulate the forward diffusion process without drift  $x_0 \sim \mathcal{N}(x_1, \sigma)$ .  
 239 After this setup, for reasonable values of  $\sigma$ , we assume that one can reverse the diffusion process and  
 240 learn the push-forward map from  $p_0(x_0 | i)$  to  $p_1(x_1 | i)$  for every index  $i$ . For this task, given the  $i$ -th  
 241 population index we denote  $p_0(x_0 | i)$  as the *source* population  $p_1(x_1 | i)$  as the  $i$ -th *target* population.

242 To construct  $p_1(x_1 | i)$ , we discretize samples from a defined silhouette; e.g. an image of a character,  
 243 where  $i$  indexes the respective character. We use upper case letters as the silhouette and generate  
 244 the corresponding samples  $x_1 \sim p_1(x_1 | i)$  from the uniform distribution over the silhouette and run  
 245 the diffusion process for samples  $x_1$  to acquire  $x_0$ . We construct the *training data* using 10 random  
 246 orientations of 24 letters, while only using the upright orientation for the remaining letters “X” and  
 247 “Y”. We construct the *test data* by using 10 random orientations of “X” and “Y” (validation and test,  
 248 respectively) that differ from the upright orientations of the same letters in the training data. We  
 249 do this to simplify the generalization task – the model will see the shapes of “X” and “Y” during  
 250 training, but the same letters under different orientations remain unseen.

251 We train FM, CGFM and 4 variants of MFM of varying  $k$  for the GCN population embedding model  
 252  $\varphi(p_t; \theta, k)$ . When  $k = 0$ ,  $\varphi(p_t; \theta, k)$  becomes identical to the DeepSets model (Zaheer et al., 2017).  
 253 We compare MFM to Flow-Matching (FM) and Conditional Generation via Flow-Matching (CGFM).  
 254 FM does not have access to conditional information; hence will only learn an aggregated lens of the  
 255 distribution dynamics and will not be able to fit the training data, and consequently won’t generalize  
 256 to the test conditions. For the training data, CGFM vector field model takes in the distribution index  
 257  $i$  as a one-hot input condition. On the test set, since none of these indices is present, we input the  
 258 normalized constant vector, which averages the learned embeddings of the indices. Because of this,  
 259 CGFM will fit the training data, however, will not be able to generalize to the unseen condition in  
 260 the test dataset. Note that the CGFM can be viewed as an *idealized* model for the train data since

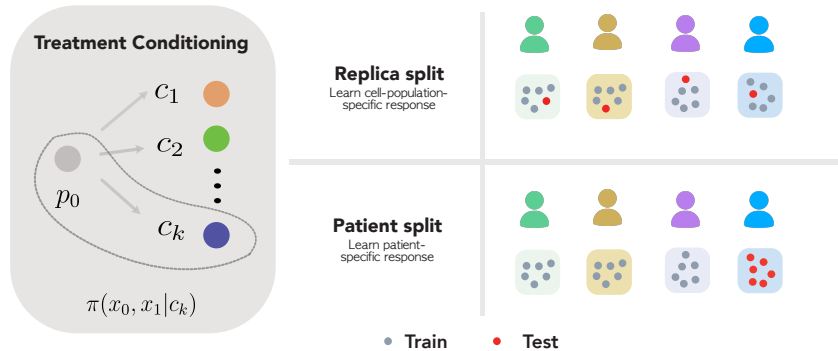


Figure 3: Organoid drug-screen dataset overview. *Left*: a given replica consists of a control distribution  $p_0$  and corresponding treatment response distribution  $p_1$  for treatment condition  $c_i$ . *Right*: train and test data splits for replica (top) and patients (bottom) splits, restively. For each experiment there are 11 treatments, 10 patients and 3 culture conditions.

261 it gets perfect information regarding the population conditions. We use CGFM to assess if other  
 262 models are fitting the data. For MFM, we expect to both fit the training data and generalize to unseen  
 263 distributional conditions.

264 In Fig. 2, we observe that indeed FM fails to adequately learn to sample from  $p_1(x_1 | i)$  in the training  
 265 set, and likewise fails to generalize, while CGFM is able to effectively sample from  $p_1(x_1 | i)$   
 266 in the training set, but fails to generalize. We report results for the synthetic experiment in Table 1.  
 267 As expected, CGFM fits the training data, however, fails to generalize beyond its set of training  
 268 conditions. In contrast, we see that MFM is able to both fit the training data (approaching the  
 269 performance of CGFM) while also generalizing to the unseen test distributions. FM fails to fit the  
 270 train data and fails to generalize under the test conditions. Interestingly, although MFM performs  
 271 better for certain values of  $k$  versus others, overall performance does not vary significantly for the  
 272 range considered.

## 273 5.2 Experiments on Organoid Drug-screen Data

274 **Data.** For experiments on biological data, we use the organoid drug-screen dataset from Ramos Zap-  
 275 atero et al. (2023). This dataset is a single-cell mass-cytometry dataset collected over 10 patients.  
 276 Somewhat unique to this dataset, unlike many prior perturbation-screen datasets which have a single  
 277 control population, this dataset has matched controls to each experimental condition. Populations from  
 278 each patient are treated with 11 different drug treatments of varying dose concentrations.<sup>1</sup> We use the  
 279 term *replicate* to define control-treatment population pairs,  $p_0(x_0 | c_i)$  and  $p_1(x_1 | c_i)$ , respectively  
 280 (see Fig. 3-left). In each patient, cell population are categorized into 3 cell *cultures*: (i) cancer associ-  
 281 ated Fibroblasts, (ii) patient-derived organoid cancer cells (PDO), and (iii) patient-derived organoid  
 282 cancer cells co-cultured fibroblasts (PDOF). We report results averaged over Fibroblast/PDO/PDOF  
 283 cultures and results for the individual cultures (this is reported in Appendix F).

284 **Pre-processing and data splits.** We filter each cell population to contain at least 1000 cells and  
 285 consider 43 bio-markers. We consider two data splits for the organoid drug-screen dataset (see  
 286 Fig. 3-right). (1) *Replicate split*; here we leave-out replicates evenly across all patients for testing. (2)  
 287 *Patients split*; here we leave-out replicates fully in one patients – in this setting, we are testing the  
 288 ability of of model to generalize population prediction of treatment response for unseen patients. In  
 289 both settings, we normalize the data and embed it into a lower dimensional principle components  
 290 (PC) representation. We do this to reduce the dimensionality of the data and to extract the relevant  
 291 information from the 43 bio-markers (features) of the ambient space. We train and evaluate all models  
 292 in the PC space. For all organoid drug-screen dataset experiments we use PC=10. Further details  
 293 regarding data pre-processing and data splits are provided in Appendix B.2.

294 For the organoid drug-screen experiments, we consider an ICNN architecture in addition to the  
 295 Flow-matching models. The ICNN model is based on CellIOT (Bunne et al., 2023a); a method for  
 296 learning cell specific response to treatments. The ICNN (and likewise CellIOT) counterparts our FM

<sup>1</sup>We consider only the highest dosage and leave exploration of dose-dependent response to future work.



Table 2: Experimental results on the organoid drug-screen dataset for population prediction of treatment response across *replicate* populations averaged over co-culture conditions. Results are reported for models trained on data embedded into 10 principle components. We report the the 1-Wasserstein ( $\mathcal{W}_1$ ), 2-Wasserstein ( $\mathcal{W}_2$ ), and the maximum-mean-discrepancy (MMD) distributional distances. We consider two settings for MFM with varying nearest-neighbours parameter. For extended results in Table 4.

	Train			Test		
	$\mathcal{W}_1$	$\mathcal{W}_2$	MMD ( $\times 10^{-3}$ )	$\mathcal{W}_1$	$\mathcal{W}_2$	MMD ( $\times 10^{-3}$ )
FM	1.946 $\pm$ 0.083	2.178 $\pm$ 0.092	6.32 $\pm$ 0.36	2.087 $\pm$ 0.035	2.301 $\pm$ 0.043	9.29 $\pm$ 0.77
ICNN	2.112 $\pm$ 0.012	2.317 $\pm$ 0.011	190.17 $\pm$ 4.87	2.200 $\pm$ 0.011	2.395 $\pm$ 0.010	249.33 $\pm$ 4.67
CGFM	<b>1.823 <math>\pm</math> 0.126</b>	<b>2.009 <math>\pm</math> 0.143</b>	<b>4.16 <math>\pm</math> 1.00</b>	2.213 $\pm$ 0.137	2.416 $\pm$ 0.154	13.91 $\pm$ 2.41
MFM ( $k = 0$ )	1.829 $\pm$ 0.050	2.012 $\pm$ 0.058	4.64 $\pm$ 0.66	1.959 $\pm$ 0.050	2.144 $\pm$ 0.059	7.35 $\pm$ 1.20
MFM ( $k = 10$ )	1.842 $\pm$ 0.049	2.020 $\pm$ 0.057	4.76 $\pm$ 0.66	<b>1.954 <math>\pm</math> 0.047</b>	<b>2.136 <math>\pm</math> 0.052</b>	<b>7.34 <math>\pm</math> 0.93</b>

Table 3: Experimental results on the organoid drug-screen dataset for population prediction of treatment response across *patient* populations. Results shown in this table are broken out in Table 5.

	Train			Test		
	$\mathcal{W}_1$	$\mathcal{W}_2$	MMD ( $\times 10^{-3}$ )	$\mathcal{W}_1$	$\mathcal{W}_2$	MMD ( $\times 10^{-3}$ )
FM	1.995 $\pm$ 0.138	2.246 $\pm$ 0.193	6.87 $\pm$ 2.65	2.607 $\pm$ 0.028	2.947 $\pm$ 0.050	21.58 $\pm$ 1.02
ICNN	2.163 $\pm$ 0.067	2.367 $\pm$ 0.070	192.67 $\pm$ 4.22	2.702 $\pm$ 0.027	2.996 $\pm$ 0.033	452.67 $\pm$ 19.14
CGFM	<b>1.773 <math>\pm</math> 0.072</b>	<b>1.954 <math>\pm</math> 0.092</b>	<b>3.03 <math>\pm</math> 0.69</b>	2.675 $\pm$ 0.019	2.938 $\pm$ 0.020	23.75 $\pm$ 0.61
MFM ( $k = 0$ )	1.863 $\pm$ 0.056	2.048 $\pm$ 0.063	5.01 $\pm$ 0.53	2.393 $\pm$ 0.160	2.685 $\pm$ 0.122	16.66 $\pm$ 1.99
MFM ( $k = 10$ )	1.881 $\pm$ 0.071	2.074 $\pm$ 0.091	5.25 $\pm$ 0.78	<b>2.326 <math>\pm</math> 0.072</b>	<b>2.610 <math>\pm</math> 0.073</b>	<b>14.30 <math>\pm</math> 2.27</b>

297 model in that it does not take the population index  $i$  as a condition. Therefore, it will neither be able  
 298 to fit the training data, nor generalize.

299 **Predicting treatment response across replicates.** We show results for generalization across repli-  
 300 cates in Table 2. As expected, we observe that CGFM fits the training data, but does not generalize to  
 301 the test replicates. With this, we can observe that the FM and ICNN models fail to fit the train data,  
 302 relative to CGFM, and also fail to generalize. MFM ( $k = 10$ ) performs best on generalization to  
 303 unseen replicates. We include results reported for the separate cell cultures in Table 4 in Appendix F.

304 **Predicting treatment response across patients.** We show results for generalization across patients  
 305 in Table 3. Similar to the replicates data setting, we observe that CGFM fits the training data, but  
 306 does not generalize to the test replicates. Likewise, the FM and ICNN models fail to fit the train data,  
 307 relative to CGFM, and also fail to generalize. MFM ( $k = 10$ ) performs best on generalization to  
 308 unseen replicates. We include results reported for the separate cell cultures in Table 5 in Appendix F.

309 Through the biological and synthetic experiments, we have shown that MFM is able to generalize  
 310 to unseen distributions/populations. The implication of our results suggest that MFM can learn  
 311 population dynamics in unseen environments. In biological contexts, like the one we have shown  
 312 in this work, this result indicates that we can learn population dynamics, of treatment response or  
 313 any arbitrary perturbation, in new/unseen patients. This works towards a model where it is possible  
 314 to predict and design an individualized treatment regimen for each patient based on their individual  
 315 characteristics and tumor microenvironment.

## 316 6 Conclusion and Future Work

317 Our paper highlights the significance of modeling dynamics based on the entire distribution. While  
 318 flow-based models offer a promising avenue for learning dynamics at the population level, they were  
 319 previously restricted to a single initial population and predefined conditions.

320 In this paper, we introduce Meta Flow Matching (MFM) as a practical solution to address these  
 321 limitations. By integrating along vector fields of the Wasserstein manifold, MFM allows for a more  
 322 comprehensive model of dynamical systems with interacting particles. Crucially, MFM leverages  
 323 graph neural networks to embed the initial population, enabling the model to generalize over various  
 324 initial distributions. MFM opens up new possibilities for understanding complex phenomena that  
 325 emerge from interacting systems in biological and physical systems.

326 In practice, we demonstrate that MFM learns meaningful embeddings of single-cell populations along  
 327 with the developmental model of these populations. Moreover, our empirical study demonstrates the  
 328 possibility of modeling patient-specific response to treatments via the meta-learning.

329 **References**

- 330 Albergo, M. S. and Vanden-Eijnden, E. (2022). Building normalizing flows with stochastic inter-  
331 polants. *arXiv preprint arXiv:2209.15571*.
- 332 Ambrosio, L., Gigli, N., and Savaré, G. (2008). *Gradient flows: in metric spaces and in the space of*  
333 *probability measures*. Springer Science & Business Media.
- 334 Amos, B., Cohen, S., Luise, G., and Redko, I. (2022). Meta optimal transport. *arXiv preprint*  
335 *arXiv:2206.05262*.
- 336 Amos, B. et al. (2023). Tutorial on amortized optimization. *Foundations and Trends® in Machine*  
337 *Learning*, 16(5):592–732.
- 338 Armingol, E., Officer, A., Harismendy, O., and Lewis, N. E. (2020). Deciphering cell–cell interactions  
339 and communication from gene expression. *Nature Reviews Genetics*, 22(2):71–88.
- 340 Benamou, J.-D. (2003). Numerical resolution of an “unbalanced” mass transport problem. *ESAIM:*  
341 *Mathematical Modelling and Numerical Analysis*, 37(5):851–868.
- 342 Binnewies, M., Roberts, E. W., Kersten, K., Chan, V., Fearon, D. F., Merad, M., Coussens, L. M.,  
343 Gabrilovich, D. I., Ostrand-Rosenberg, S., Hedrick, C. C., Vonderheide, R. H., Pittet, M. J., Jain,  
344 R. K., Zou, W., Howcroft, T. K., Woodhouse, E. C., Weinberg, R. A., and Krummel, M. F. (2018).  
345 Understanding the tumor immune microenvironment (time) for effective therapy. *Nature Medicine*,  
346 24(5):541–550.
- 347 Bunne, C., Stark, S. G., Gut, G., Del Castillo, J. S., Levesque, M., Lehmann, K.-V., Pelkmans, L.,  
348 Krause, A., and Räscher, G. (2023a). Learning single-cell perturbation responses using neural  
349 optimal transport. *Nature Methods*, 20(11):1759–1768.
- 350 Bunne, C., Stark, S. G., Gut, G., del Castillo, J. S., Levesque, M., Lehmann, K.-V., Pelkmans, L.,  
351 Krause, A., and Räscher, G. (2023b). Learning single-cell perturbation responses using neural  
352 optimal transport. *Nature Methods*, 20(11):1759–1768.
- 353 Chen, S., Rivaud, P., Park, J. H., Tsou, T., Charles, E., Haliburton, J. R., Pichiorri, F., and Thomson,  
354 M. (2020). Dissecting heterogeneous cell populations across drug and disease conditions with  
355 popalign. *Proceedings of the National Academy of Sciences*, 117(46):28784–28794.
- 356 Chen, T., Chen, X., Chen, W., Heaton, H., Liu, J., Wang, Z., and Yin, W. (2022). Learning to optimize:  
357 A primer and a benchmark. *Journal of Machine Learning Research*, 23(189):1–59.
- 358 Chizat, L., Peyré, G., Schmitzer, B., and Vialard, F.-X. (2018). Unbalanced optimal transport:  
359 Dynamic and kantorovich formulations. *Journal of Functional Analysis*, 274(11):3090–3123.
- 360 Chung, W., Eum, H. H., Lee, H.-O., Lee, K.-M., Lee, H.-B., Kim, K.-T., Ryu, H. S., Kim, S., Lee, J. E.,  
361 Park, Y. H., Kan, Z., Han, W., and Park, W.-Y. (2017). Single-cell rna-seq enables comprehensive  
362 tumour and immune cell profiling in primary breast cancer. *Nature Communications*, 8(1).
- 363 Dao, Q., Phung, H., Nguyen, B., and Tran, A. (2023). Flow matching in latent space. *arXiv preprint*  
364 *arXiv:2307.08698*.
- 365 De Bortoli, V., Thornton, J., Heng, J., and Doucet, A. (2021). Diffusion schrödinger bridge with  
366 applications to score-based generative modeling. *Advances in Neural Information Processing*  
367 *Systems*, 34:17695–17709.
- 368 Frangieh, C. J., Melms, J. C., Thakore, P. I., Geiger-Schuller, K. R., Ho, P., Luoma, A. M., Cleary, B.,  
369 Jerby-Aron, L., Malu, S., Cuoco, M. S., Zhao, M., Ager, C. R., Rogava, M., Hovey, L., Rotem,  
370 A., Bernatchez, C., Wucherpfennig, K. W., Johnson, B. E., Rozenblatt-Rosen, O., Schadendorf,  
371 D., Regev, A., and Izar, B. (2021). Multimodal pooled perturb-cite-seq screens in patient models  
372 define mechanisms of cancer immune evasion. *Nature Genetics*, 53(3):332–341.
- 373 Gilmer, J., Schoenholz, S. S., Riley, P. F., Vinyals, O., and Dahl, G. E. (2017). Neural message  
374 passing for quantum chemistry. In *International conference on machine learning*, pages 1263–1272.  
375 PMLR.

376 Goodenough, D. A. and Paul, D. L. (2009). Gap junctions. *Cold Spring Harb Perspect Biol*,  
377 1(1):a002576.

378 Gretton, A., Borgwardt, K. M., Rasch, M. J., Schölkopf, B., and Smola, A. (2012). A kernel  
379 two-sample test. *The Journal of Machine Learning Research*, 13(1):723–773.

380 Gulati, G. S., Sikandar, S. S., Wesche, D. J., Manjunath, A., Bharadwaj, A., Berger, M. J., Ilagan, F.,  
381 Kuo, A. H., Hsieh, R. W., Cai, S., Zabala, M., Scheeren, F. A., Lobo, N. A., Qian, D., Yu, F. B.,  
382 Dirbas, F. M., Clarke, M. F., and Newman, A. M. (2020). Single-cell transcriptional diversity is a  
383 hallmark of developmental potential. *Science*, 367(6476):405–411.

384 Hashimoto, T. B., Gifford, D. K., and Jaakkola, T. S. (2016). Learning population-level diffusions  
385 with generative recurrent networks. In *Proceedings of the 33rd International Conference on*  
386 *Machine Learning*, pages 2417–2426.

387 Hetzel, L., Boehm, S., Kilbertus, N., Günemann, S., Lotfollahi, M., and Theis, F. (2022). Predicting  
388 cellular responses to novel drug perturbations at a single-cell resolution. In Koyejo, S., Mohamed,  
389 S., Agarwal, A., Belgrave, D., Cho, K., and Oh, A., editors, *Advances in Neural Information*  
390 *Processing Systems*, volume 35, pages 26711–26722. Curran Associates, Inc.

391 Huguet, G., Magruder, D. S., Tong, A., Fasina, O., Kuchroo, M., Wolf, G., and Krishnaswamy, S.  
392 (2022). Manifold interpolating optimal-transport flows for trajectory inference.

393 Huguet, G., Tong, A., Zapatero, M. R., Wolf, G., and Krishnaswamy, S. (2023). Geodesic sinkhorn:  
394 Optimal transport for high-dimensional datasets. In *IEEE MLSP*.

395 Isobe, N., Koyama, M., Hayashi, K., and Fukumizu, K. (2024). Extended flow matching: a method  
396 of conditional generation with generalized continuity equation. *arXiv preprint arXiv:2402.18839*.

397 Ji, Y., Lotfollahi, M., Wolf, F. A., and Theis, F. J. (2021). Machine learning for perturbational  
398 single-cell omics. *Cell Systems*, 12(6):522–537.

399 Koshizuka, T. and Sato, I. (2023). Neural lagrangian schrödinger bridge. In *ICLR*.

400 Lipman, Y., Chen, R. T. Q., Ben-Hamu, H., Nickel, M., and Le, M. (2023). Flow matching for  
401 generative modeling. In *The Eleventh International Conference on Learning Representations*.

402 Liu, G.-H., Vahdat, A., Huang, D.-A., Theodorou, E. A., Nie, W., and Anandkumar, A. (2023).  $I^2sb$ :  
403 Image-to-image schrödinger bridge. In *ICML*.

404 Liu, X., Gong, C., and Liu, Q. (2022). Flow straight and fast: Learning to generate and transfer data  
405 with rectified flow. *arXiv preprint arXiv:2209.03003*.

406 Lopez, R., Regier, J., Cole, M. B., Jordan, M. I., and Yosef, N. (2018). Deep generative modeling for  
407 single-cell transcriptomics. *Nature Methods*, 15(12):1053–1058.

408 Lotfollahi, M., Wolf, F. A., and Theis, F. J. (2019). scgen predicts single-cell perturbation responses.  
409 *Nature Methods*, 16(8):715–721.

410 Makkuva, A. V., Taghvaei, A., Oh, S., and Lee, J. D. (2020). Optimal transport mapping via input  
411 convex neural networks. In *ICML*.

412 Molè, M. A., Coorens, T. H. H., Shahbazi, M. N., Weberling, A., Weatherbee, B. A. T., Gantner,  
413 C. W., Sancho-Serra, C., Richardson, L., Drinkwater, A., Syed, N., Engley, S., Snell, P., Christie,  
414 L., Elder, K., Campbell, A., Fishel, S., Behjati, S., Vento-Tormo, R., and Zernicka-Goetz, M.  
415 (2021). A single cell characterisation of human embryogenesis identifies pluripotency transitions  
416 and putative anterior hypoblast centre. *Nature Communications*, 12(1).

417 Neklyudov, K., Brekelmans, R., Tong, A., Atanackovic, L., Liu, Q., and Makhzani, A. (2023). A com-  
418 putational framework for solving wasserstein lagrangian flows. *arXiv preprint arXiv:2310.10649*.

419 Neklyudov, K., Severo, D., and Makhzani, A. (2022). Action matching: A variational method for  
420 learning stochastic dynamics from samples.

421 Otto, F. (2001). The geometry of dissipative evolution equations: the porous medium equation.

- 422 Peidli, S., Green, T. D., Shen, C., Gross, T., Min, J., Garda, S., Yuan, B., Schumacher, L. J., Taylor-  
423 King, J. P., Marks, D. S., et al. (2024). scperturb: harmonized single-cell perturbation data. *Nature*  
424 *Methods*, pages 1–10.
- 425 Peyré, G. and Cuturi, M. (2019). *Computational Optimal Transport*. arXiv:1803.00567.
- 426 Pooladian, A.-A., Ben-Hamu, H., Domingo-Enrich, C., Amos, B., Lipman, Y., and Chen, R. T.  
427 (2023). Multisample flow matching: Straightening flows with minibatch couplings. *arXiv preprint*  
428 *arXiv:2304.14772*.
- 429 Price, I., Sanchez-Gonzalez, A., Alet, F., Ewalds, T., El-Kadi, A., Stott, J., Mohamed, S., Battaglia,  
430 P., Lam, R., and Willson, M. (2023). Gencast: Diffusion-based ensemble forecasting for medium-  
431 range weather. *arXiv preprint arXiv:2312.15796*.
- 432 Ramos Zapatero, M., Tong, A., Opzooomer, J. W., O’Sullivan, R., Cardoso Rodriguez, F., Sufi, J.,  
433 Vlckova, P., Nattress, C., Qin, X., Claus, J., Hochhauser, D., Krishnaswamy, S., and Tape, C. J.  
434 (2023). Trellis tree-based analysis reveals stromal regulation of patient-derived organoid drug  
435 responses. *Cell*, 186(25):5606–5619.e24.
- 436 Rizvi, A. H., Camara, P. G., Kandror, E. K., Roberts, T. J., Schieren, I., Maniatis, T., and Rabadan, R.  
437 (2017). Single-cell topological rna-seq analysis reveals insights into cellular differentiation and  
438 development. *Nature Biotechnology*, 35(6):551–560.
- 439 Rombach, R., Blattmann, A., Lorenz, D., Esser, P., and Ommer, B. (2022). High-resolution image  
440 synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer*  
441 *vision and pattern recognition*, pages 10684–10695.
- 442 Saharia, C., Chan, W., Chang, H., Lee, C., Ho, J., Salimans, T., Fleet, D., and Norouzi, M. (2022a).  
443 Palette: Image-to-image diffusion models. In *ACM SIGGRAPH 2022 conference proceedings*,  
444 pages 1–10.
- 445 Saharia, C., Chan, W., Saxena, S., Li, L., Whang, J., Denton, E. L., Ghasemipour, K., Gontijo Lopes,  
446 R., Karagol Ayan, B., Salimans, T., et al. (2022b). Photorealistic text-to-image diffusion models  
447 with deep language understanding. *Advances in neural information processing systems*, 35:36479–  
448 36494.
- 449 Sanchez-Gonzalez, A., Godwin, J., Pfaff, T., Ying, R., Leskovec, J., and Battaglia, P. (2020). Learning  
450 to simulate complex physics with graph networks. In *International conference on machine learning*,  
451 pages 8459–8468. PMLR.
- 452 Schiebinger, G., Shu, J., Tabaka, M., Cleary, B., Subramanian, V., Solomon, A., Gould, J., Liu,  
453 S., Lin, S., Berube, P., et al. (2019). Optimal-transport analysis of single-cell gene expression  
454 identifies developmental trajectories in reprogramming. *Cell*, 176(4):928–943.
- 455 Somnath, V. R., Pariset, M., Hsieh, Y.-P., Martinez, M. R., Krause, A., and Bunne, C. (2023). Aligned  
456 diffusion schr"odinger bridges. In *UAI*.
- 457 Tong, A., FATRAS, K., Malkin, N., Huguét, G., Zhang, Y., Rector-Brooks, J., Wolf, G., and Bengio,  
458 Y. (2024). Improving and generalizing flow-based generative models with minibatch optimal  
459 transport. *Transactions on Machine Learning Research*. Expert Certification.
- 460 Tong, A., Huang, J., Wolf, G., Van Dijk, D., and Krishnaswamy, S. (2020). Trajectorynet: A dynamic  
461 optimal transport network for modeling cellular dynamics. In *International conference on machine*  
462 *learning*, pages 9526–9536. PMLR.
- 463 Tong, A., Malkin, N., Fatras, K., Atanackovic, L., Zhang, Y., Huguét, G., Wolf, G., and Bengio,  
464 Y. (2023). Simulation-free schr"odinger bridges via score and flow matching. *arXiv preprint*  
465 *arXiv:2307.03672*.
- 466 Uscidda, T. and Cuturi, M. (2023). The monge gap: A regularizer to learn all transport maps. In  
467 Krause, A., Brunskill, E., Cho, K., Engelhardt, B., Sabato, S., and Scarlett, J., editors, *Proceedings*  
468 *of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine*  
469 *Learning Research*, pages 34709–34733. PMLR.

- 470 Verma, Y., Heinonen, M., and Garg, V. (2024). Climode: Climate and weather forecasting with  
471 physics-informed neural odes. *arXiv preprint arXiv:2404.10024*.
- 472 Villani, C. (2009). *Optimal transport: old and new*, volume 338. Springer.
- 473 Weinreb, C., Wolock, S., Tusi, B. K., Socolovsky, M., and Klein, A. M. (2018). Fundamental limits  
474 on dynamic inference from single-cell snapshots. 115(10):E2467–E2476.
- 475 Yang, K. D. and Uhler, C. (2019). Scalable unbalanced optimal transport using generative adversarial  
476 networks. In *7th International Conference on Learning Representations*, page 20.
- 477 Zaheer, M., Kottur, S., Ravanbakhsh, S., Poczos, B., Salakhutdinov, R. R., and Smola, A. J. (2017).  
478 Deep sets. *Advances in neural information processing systems*, 30.
- 479 Zeng, T. and Dai, H. (2019). Single-cell rna sequencing-based computational analysis to describe  
480 disease heterogeneity. *Frontiers in Genetics*, 10.
- 481 Zheng, Q., Le, M., Shaul, N., Lipman, Y., Grover, A., and Chen, R. T. (2023). Guided flows for  
482 generative modeling and decision making. *arXiv preprint arXiv:2311.13443*.

## 483 A Proof of Theorem 1

484 **Theorem 1.** Consider a dataset of populations  $\mathcal{D} = \{(\pi(x_0, x_1 | i))\}_i$  generated from some unknown  
 485 conditional model  $\pi(x_0, x_1 | c)p(c)$ . Then the following objective

$$\mathcal{L}(\omega, \theta) = \mathbb{E}_{p(c)} \int_0^1 dt \mathbb{E}_{p_t(x_t | c)} \|v_t^*(x_t | c) - v_t(x_t | \varphi(p_0, \theta), \omega)\|^2 \quad (16)$$

486 is equivalent to the Meta Flow Matching objective

$$\mathcal{L}_{MFM}(\omega, \theta) = \mathbb{E}_{i \sim \mathcal{D}} \mathbb{E}_{\pi(x_0, x_1 | i)} \int_0^1 dt \left\| \frac{\partial}{\partial t} f_t(x_0, x_1) - v_t(f_t(x_0, x_1) | \varphi(p_0, \theta); \omega) \right\|^2 \quad (17)$$

487 up to an additive constant.

488 *Proof.* The loss function

$$\mathcal{L}(\omega, \theta) = \mathbb{E}_{p(c)} \int_0^1 dt \mathbb{E}_{p_t(x_t | c)} \|v_t^*(x_t | c) - v_t(x_t | \varphi(p_t, \theta); \omega)\|^2 \quad (20)$$

$$= -2\mathbb{E}_{p(c)} \int dt dx \langle p_t(x | c) v_t^*(x | c), v_t(x | \varphi(p_t, \theta); \omega) \rangle + \quad (21)$$

$$+ \mathbb{E}_{p(c)} \int_0^1 dt \mathbb{E}_{p_t(x_t | c)} \|v_t(x_t | \varphi(p_t, \theta), \omega)\|^2 + \quad (22)$$

$$+ \mathbb{E}_{p(c)} \int_0^1 dt \mathbb{E}_{p_t(x_t | c)} \|v_t^*(x_t | c)\|^2. \quad (23)$$

489 The last term does not depend on  $\theta$ , the second term we can estimate, for the first term, we use the  
 490 formula for the (from Eq. (8))

$$p_t(\xi | c) v_t^*(\xi | c) = \mathbb{E}_{\pi(x_0, x_1)} \delta(f_t(x_0, x_1) - \xi) \frac{\partial f_t(x_0, x_1)}{\partial t}. \quad (24)$$

491 Thus, the loss is equivalent (up to a constant) to

$$\mathcal{L}(\omega, \theta) = -2\mathbb{E}_{p(c)} \mathbb{E}_{\pi(x_0, x_1 | c)} \int dt \left\langle \frac{\partial f_t(x_0, x_1)}{\partial t}, v_t(f_t(x_0, x_1) | \varphi(p_t, \theta); \omega) \right\rangle + \quad (25)$$

$$+ \mathbb{E}_{p(c)} \mathbb{E}_{\pi(x_0, x_1 | c)} \int_0^1 dt \|v_t(f_t(x_0, x_1) | \varphi(p_t, \theta), \omega)\|^2 \pm \quad (26)$$

$$\pm \mathbb{E}_{p(c)} \mathbb{E}_{\pi(x_0, x_1 | c)} \int_0^1 dt \left\| \frac{\partial f_t(x_0, x_1)}{\partial t} \right\|^2 \quad (27)$$

$$= \mathbb{E}_{c \sim p(c)} \mathbb{E}_{\pi(x_0, x_1 | c)} \int_0^1 dt \left\| \frac{\partial}{\partial t} f_t(x_0, x_1) - v_t(f_t(x_0, x_1) | \varphi(p_t, \theta); \omega) \right\|^2. \quad (28)$$

492 Note that in the final expression we do not need access to the probabilistic model of  $p(c)$  if the joints  
 493  $\pi(x_0, x_1 | c)$  are already sampled in the data  $\mathcal{D}$ . Thus, we have

$$\mathcal{L}(\omega, \theta) = \mathbb{E}_{c \sim p(c)} \mathbb{E}_{\pi(x_0, x_1 | c)} \int_0^1 dt \left\| \frac{\partial}{\partial t} f_t(x_0, x_1) - v_t(f_t(x_0, x_1) | \varphi(p_t, \theta); \omega) \right\|^2 \quad (29)$$

$$= \mathbb{E}_{i \sim \mathcal{D}} \mathbb{E}_{\pi(x_0, x_1 | i)} \int_0^1 dt \left\| \frac{\partial}{\partial t} f_t(x_0, x_1) - v_t(f_t(x_0, x_1) | \varphi(p_t, \theta); \omega) \right\|^2 \quad (30)$$

$$= \mathcal{L}_{MFM}(\omega, \theta). \quad (31)$$

494 □

## 495 B Experimental Details

### 496 B.1 Synthetic letters data

497 The synthetic letters dataset contains 242 train populations a 10 test populations. Each population  
 498 contains roughly between 750 and 2700 samples. In this dataset.

## 499 B.2 Organoid drug-screen data

500 The organoid drug-screen dataset contains a total of 927 replicates (or coupled populations). In the  
501 *replicates split*, we use 713 populations for training and 103 left-out populations for testing. In the  
502 *patients split*, we use 861 populations for training and 33 left-out populations for testing.

## 503 B.3 Model architectures and hyperparameters

504 **ICNN.** The ICNN baseline was constructed with two networks ICNN network  $f(x)$  and  $g(x)$ , with  
505 non-negative leaky ReLU activation layers.  $f(x)$  is used to minimize the transport distance and  $g(x)$   
506 is used to transport from source to target. It has four hidden units with width of 64, and a latent  
507 dimension of 50. Both networks uses Adam optimizer (lr=1e-4,  $\beta_1=0.5$ ,  $\beta_2=0.9$ ).  $g(x)$  is trained  
508 with an inner iteration of 10 for every iteration  $f(x)$  is trained.

509 **Vector Field Models.** All vector field models  $v_t$  are parameterized 4 linear layers with 512 hidden  
510 units and SELU activation functions. The FM vector field model additionally takes a conditional  
511 input for the one-hot treatment encoding. CGFM takes the conditional input for the one-hot treatment  
512 conditions as well as a one-hot encoding for the population index condition  $i$ . The MFM vector field  
513 model takes population embedding conditions, that is output from the GCN, as input, as well as the  
514 treatment one-hot encoding. All vector field models use temporal embeddings for time and positional  
515 embeddings for the input samples. We did not sweep the size of this embeddings space and found  
516 that a temporal embedding and positional embeddings sizes of 128 worked sufficiently well.

517 **Graph Neural Network.** We considered a GCN model that consists of a  $k$ -nearest neighbour graph  
518 edge pooling layer and 3 graph convolution layers with 512 hidden units. The final GCN model  
519 layer outputs an embedding representation  $e \in \mathbb{R}^d$ . For the Synthetic experiment, we found that  
520  $d = 256$  performed well, and  $d = 128$  performed well for the biological experiments. We normalize  
521 and project embeddings onto a hyper-sphere, and find that this normalization helps improve training.  
522 Additionally, the GCN takes a one-hot cell-type encoding (encoding for Fibroblast cells or PDO  
523 cells) for the control populations  $p_0$ . This may be beneficial for PDOF populations where both  
524 Fibroblast cells and PDO cells are present. However, it is important to note that labeling which cells  
525 are Fibroblasts versus PDOs within the PDOF cultures is difficult and noisy in itself, hence such a  
526 cell-type condition may yield no additive information/performance gain.

527 **Optimization.** We use the Adam optimizer with a learning rate of 0.0001 for all Flow-matching  
528 models (FM, CGFM, MFM). We also used the Adam optimizer with a learning rate of 0.0001 for  
529 the GCN model. To train the MFM (FM+GCN) models, we alternate between updating the vector  
530 field model parameters  $\omega$  and the GCN model parameters  $\theta$ . We alternate between updating the  
531 respective model parameters every epoch. FM and CGFM model were trained for 2000 epochs, while  
532 MFM models were trained for 4000 epochs. Due to the alternating optimization, the MFM vector  
533 field model receives half as many updates compared to its counterparts (FM and CGFM). Therefore,  
534 training for the double the epochs is necessary for fair comparison.

535 The hyperparameters stated in this section were selected from brief and small grid search sweeps. We  
536 did not conduct any thorough hyperparameter optimization.

## 537 C Implementation Details

538 We implement all our experiments using PyTorch and PyTorch Geometric. We submitted our code as  
539 supplementary material with our submission.

540 All experiments were conducted on a HPC cluster primarily on NVIDIA Tesla T4 16GB GPUs. Each  
541 individual seed experiment run required only 1 GPU. Each experiment ran between 3-11 hours and  
542 all experiments took approximately 500 GPU hours.

## 543 D Limitations

544 In this work we explored empirically the effect of conditioning the learned flow on the initial  
545 distribution. We argue this is a more natural model for many biological systems. However, there  
546 are many other aspects of modeling biological systems that we did not consider. In particular we

547 did not consider extensions to the manifold setting (Huguet et al., 2022, 2023), unbalanced optimal  
 548 transport (Benamou, 2003; Yang and Uhler, 2019; Chizat et al., 2018), aligned (Somnath et al., 2023;  
 549 Liu et al., 2023), or stochastic settings (Bunne et al., 2023a; Koshizuka and Sato, 2023) in this work.

550 **E Broader Impacts**

551 This paper is primarily a theoretical and methodological contribution with little societal impact. MFM  
 552 can be used to better model dynamical systems of interacting particles and in particular cellular  
 553 systems. Better modeling of cellular systems can potentially be used for the development of malicious  
 554 biological agents. However, we do not see this as a significant risk at this time.

555 **F Extended Results**

Table 4: Experimental results on the organoid drug-screen dataset for population prediction of treatment response across replicate populations. Results are reported for models trained on data embedded into 10 principle components. We report the the 1-Wasserstein ( $\mathcal{W}_1$ ), 2-Wasserstein ( $\mathcal{W}_2$ ), and the maximum-mean-discrepancy (MMD) distributional distances. We consider 2 settings for MFM with varying nearest-neighbours parameter.

<b>Fibroblasts</b>						
	$\mathcal{W}_1$	Train $\mathcal{W}_2$	MMD ( $\times 10^{-3}$ )	$\mathcal{W}_1$	Test $\mathcal{W}_2$	MMD ( $\times 10^{-3}$ )
FM	1.584 $\pm$ 0.022	1.730 $\pm$ 0.015	3.12 $\pm$ 0.59	1.612 $\pm$ 0.014	1.736 $\pm$ 0.024	3.62 $\pm$ 0.15
ICNN	1.613 $\pm$ 0.010	1.703 $\pm$ 0.010	52.4 $\pm$ 1.64	1.655 $\pm$ 0.008	1.746 $\pm$ 0.008	53.0 $\pm$ 5.00
CGFM	<b>1.472 <math>\pm</math> 0.046</b>	<b>1.548 <math>\pm</math> 0.048</b>	<b>1.28 <math>\pm</math> 0.74</b>	1.633 $\pm$ 0.022	1.724 $\pm$ 0.023	4.95 $\pm$ 0.72
MFM ( $k = 0$ )	1.519 $\pm$ 0.034	1.599 $\pm$ 0.036	2.56 $\pm$ 0.56	<b>1.574 <math>\pm</math> 0.002</b>	<b>1.657 <math>\pm</math> 0.003</b>	<b>3.31 <math>\pm</math> 0.12</b>
MFM ( $k = 10$ )	1.547 $\pm$ 0.027	1.617 $\pm$ 0.027	2.84 $\pm$ 0.56	1.576 $\pm$ 0.017	1.658 $\pm$ 0.019	3.44 $\pm$ 0.19
<b>PDO</b>						
	$\mathcal{W}_1$	Train $\mathcal{W}_2$	MMD ( $\times 10^{-3}$ )	$\mathcal{W}_1$	Test $\mathcal{W}_2$	MMD ( $\times 10^{-3}$ )
FM	2.002 $\pm$ 0.027	2.201 $\pm$ 0.025	6.40 $\pm$ 0.10	2.033 $\pm$ 0.015	2.210 $\pm$ 0.016	6.92 $\pm$ 0.65
ICNN	2.29 $\pm$ 0.005	2.458 $\pm$ 0.003	245.8 $\pm$ 9.18	2.247 $\pm$ 0.005	2.415 $\pm$ 0.004	153 $\pm$ 1.00
CGFM	1.818 $\pm$ 0.198	1.931 $\pm$ 0.229	3.78 $\pm$ 0.27	2.255 $\pm$ 0.216	2.434 $\pm$ 0.240	12.16 $\pm$ 3.87
MFM ( $k = 0$ )	1.817 $\pm$ 0.043	1.935 $\pm$ 0.040	<b>3.61 <math>\pm</math> 0.50</b>	1.909 $\pm$ 0.076	2.057 $\pm$ 0.098	<b>5.14 <math>\pm</math> 0.92</b>
MFM ( $k = 10$ )	<b>1.805 <math>\pm</math> 0.074</b>	<b>1.921 <math>\pm</math> 0.078</b>	3.68 $\pm$ 0.78	<b>1.903 <math>\pm</math> 0.068</b>	<b>2.051 <math>\pm</math> 0.084</b>	<b>5.14 <math>\pm</math> 0.90</b>
<b>PDOF</b>						
	$\mathcal{W}_1$	Train $\mathcal{W}_2$	MMD ( $\times 10^{-3}$ )	$\mathcal{W}_1$	Test $\mathcal{W}_2$	MMD ( $\times 10^{-3}$ )
FM	2.252 $\pm$ 0.20	2.603 $\pm$ 0.236	9.43 $\pm$ 0.38	2.616 $\pm$ 0.076	2.958 $\pm$ 0.089	19.34 $\pm$ 1.51
ICNN	2.432 $\pm$ 0.021	2.791 $\pm$ 0.020	272.3 $\pm$ 3.80	2.699 $\pm$ 0.021	3.023 $\pm$ 0.019	542 $\pm$ 8.00
CGFM	2.179 $\pm$ 0.133	2.548 $\pm$ 0.153	<b>7.42 <math>\pm</math> 2.00</b>	2.750 $\pm$ 0.173	3.089 $\pm$ 0.200	22.63 $\pm$ 2.64
MFM ( $k = 0$ )	<b>2.150 <math>\pm</math> 0.073</b>	<b>2.502 <math>\pm</math> 0.099</b>	7.75 $\pm$ 0.93	2.395 $\pm$ 0.071	2.717 $\pm$ 0.076	13.61 $\pm$ 2.56
MFM ( $k = 10$ )	2.174 $\pm$ 0.046	2.523 $\pm$ 0.067	7.75 $\pm$ 0.65	<b>2.382 <math>\pm</math> 0.055</b>	<b>2.699 <math>\pm</math> 0.054</b>	<b>13.45 <math>\pm</math> 1.69</b>



Table 5: Experimental results on the organoid drug-screen dataset for population prediction of treatment response across **patient** populations. Results are reported for models trained on data embedded into 10 principle components. We report the the 1-Wasserstein ( $\mathcal{W}_1$ ), 2-Wasserstein ( $\mathcal{W}_2$ ), and the maximum-mean-discrepancy (MMD) distributional distances. We consider 2 settings for MFM with varying nearest-neighbours parameter.

Fibroblasts						
	$\mathcal{W}_1$	Train $\mathcal{W}_2$	MMD ( $\times 10^{-3}$ )	$\mathcal{W}_1$	Test $\mathcal{W}_2$	MMD ( $\times 10^{-3}$ )
FM	$1.599 \pm 0.071$	$1.761 \pm 0.137$	$2.82 \pm 0.34$	$1.667 \pm 0.003$	$1.846 \pm 0.064$	$7.85 \pm 0.15$
ICNN	$1.695 \pm 0.08$	$1.796 \pm 0.09$	$48.2 \pm 3.412$	$1.6 \pm 0.009$	$1.68 \pm 0.013$	$62.2 \pm 1.32$
CGFM	<b><math>1.496 \pm 0.019</math></b>	<b><math>1.572 \pm 0.016</math></b>	<b><math>1.45 \pm 0.14</math></b>	$1.566 \pm 0.028$	$1.652 \pm 0.026$	$6.46 \pm 0.82$
MFM ( $k = 0$ )	$1.551 \pm 0.037$	$1.632 \pm 0.042$	$2.31 \pm 0.71$	$1.453 \pm 0.200$	$1.527 \pm 0.022$	$3.66 \pm 0.67$
MFM ( $k = 10$ )	$1.555 \pm 0.034$	$1.635 \pm 0.039$	$2.54 \pm 0.42$	<b><math>1.441 \pm 0.003</math></b>	<b><math>1.514 \pm 0.001</math></b>	<b><math>3.37 \pm 0.72</math></b>
PDO						
	$\mathcal{W}_1$	Train $\mathcal{W}_2$	MMD ( $\times 10^{-3}$ )	$\mathcal{W}_1$	Test $\mathcal{W}_2$	MMD ( $\times 10^{-3}$ )
FM	$1.996 \pm 0.196$	$2.171 \pm 0.243$	$6.79 \pm 3.40$	$2.128 \pm 0.064$	$2.312 \pm 0.075$	$7.88 \pm 1.26$
ICNN	$2.315 \pm 0.060$	$2.478 \pm 0.057$	$236.8 \pm 0.006$	$2.538 \pm 0.018$	$2.731 \pm 0.027$	$232.8 \pm 20.6$
CGFM	<b><math>1.662 \pm 0.026</math></b>	<b><math>1.760 \pm 0.023</math></b>	<b><math>1.74 \pm 0.16</math></b>	$2.460 \pm 0.018$	$2.533 \pm 0.023$	$13.6 \pm 0.25$
MFM ( $k = 0$ )	$1.837 \pm 0.058$	$1.964 \pm 0.059$	$3.74 \pm 0.29$	$2.010 \pm 0.142$	$2.168 \pm 0.182$	$6.01 \pm 1.77$
MFM ( $k = 10$ )	$1.838 \pm 0.035$	$1.957 \pm 0.038$	$3.75 \pm 0.41$	<b><math>1.971 \pm 0.082</math></b>	<b><math>2.114 \pm 0.101</math></b>	<b><math>5.42 \pm 1.11</math></b>
PDOF						
	$\mathcal{W}_1$	Train $\mathcal{W}_2$	MMD ( $\times 10^{-3}$ )	$\mathcal{W}_1$	Test $\mathcal{W}_2$	MMD ( $\times 10^{-3}$ )
FM	$2.390 \pm 0.148$	$2.806 \pm 0.198$	$11.0 \pm 2.21$	$4.026 \pm 0.018$	$4.683 \pm 0.011$	$49.0 \pm 1.66$
ICNN	$2.479 \pm 0.06$	$2.826 \pm 0.063$	$291 \pm 9.24$	$3.968 \pm 0.0554$	$4.579 \pm 0.060$	$1263 \pm 37.5$
CGFM	<b><math>2.160 \pm 0.170</math></b>	<b><math>2.530 \pm 0.237</math></b>	<b><math>7.90 \pm 1.79</math></b>	$4.000 \pm 0.010$	$4.629 \pm 0.012$	$49.2 \pm 0.76$
MFM ( $k = 0$ )	$2.202 \pm 0.072$	$2.548 \pm 0.089$	$8.98 \pm 0.59$	$3.717 \pm 0.138$	$4.360 \pm 0.162$	$40.3 \pm 3.52$
MFM ( $k = 10$ )	$2.251 \pm 0.143$	$2.631 \pm 0.197$	$9.45 \pm 1.52$	<b><math>3.565 \pm 0.132</math></b>	<b><math>4.201 \pm 0.119</math></b>	<b><math>36.1 \pm 4.97</math></b>

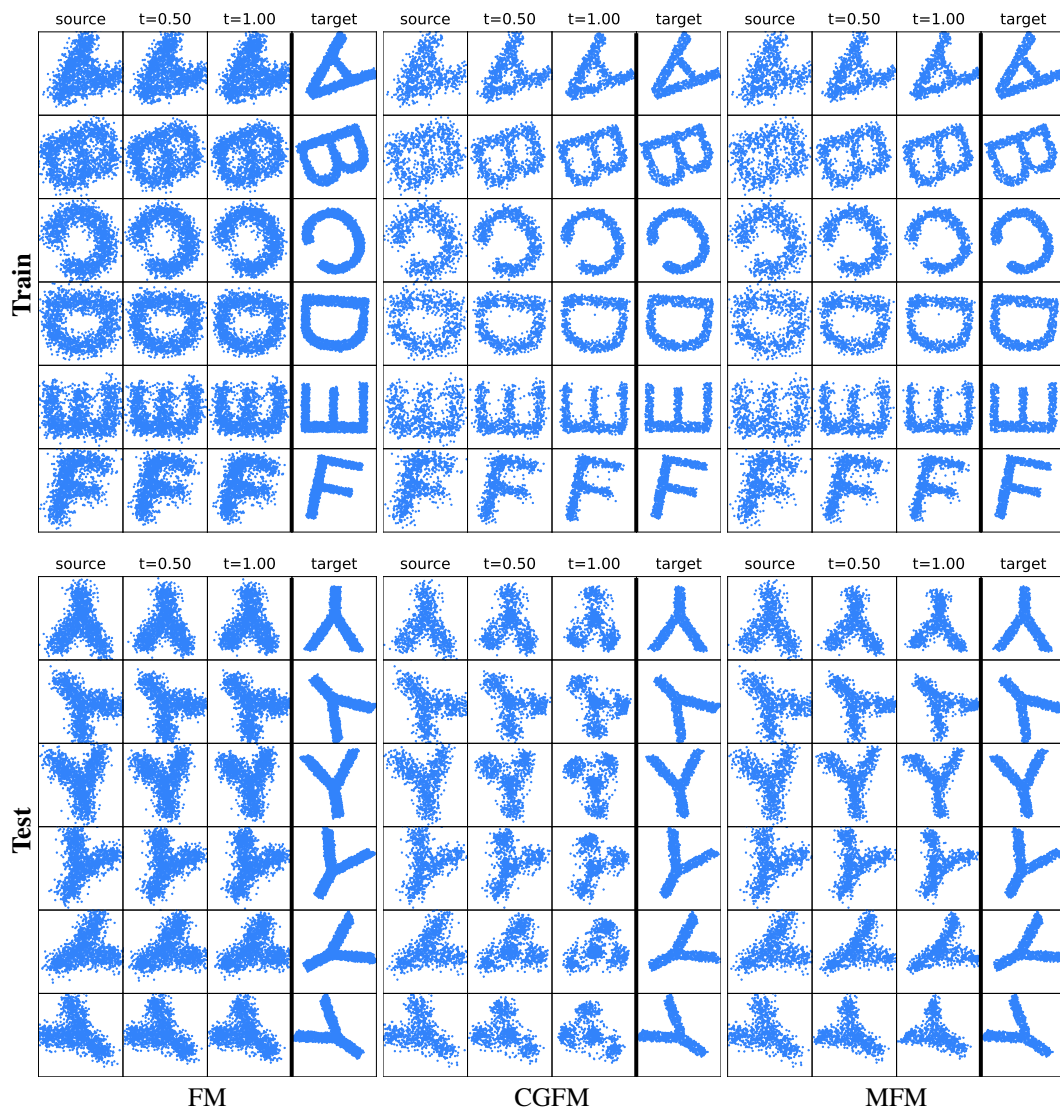


Figure 4: Model-generated samples for synthetic letters from the source ( $t = 0$ ) to target ( $t = 1$ ) distributions.

556 **NeurIPS Paper Checklist**

557 **1. Claims**

558 Question: Do the main claims made in the abstract and introduction accurately reflect the  
559 paper's contributions and scope?

560 Answer: [Yes]

561 Justification: Claims and contributions introduced in abstract and introduction are sup-  
562 ported with theoretical result in Section 3 and empirical results through synthetic and real  
563 experiments in Section 5.

564 Guidelines:

- 565 • The answer NA means that the abstract and introduction do not include the claims  
566 made in the paper.
- 567 • The abstract and/or introduction should clearly state the claims made, including the  
568 contributions made in the paper and important assumptions and limitations. A No or  
569 NA answer to this question will not be perceived well by the reviewers.
- 570 • The claims made should match theoretical and experimental results, and reflect how  
571 much the results can be expected to generalize to other settings.
- 572 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
573 are not attained by the paper.

574 **2. Limitations**

575 Question: Does the paper discuss the limitations of the work performed by the authors?

576 Answer: [Yes]

577 Justification: We discuss limitations in Appendix D.

578 Guidelines:

- 579 • The answer NA means that the paper has no limitation while the answer No means that  
580 the paper has limitations, but those are not discussed in the paper.
- 581 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 582 • The paper should point out any strong assumptions and how robust the results are to  
583 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
584 model well-specification, asymptotic approximations only holding locally). The authors  
585 should reflect on how these assumptions might be violated in practice and what the  
586 implications would be.
- 587 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
588 only tested on a few datasets or with a few runs. In general, empirical results often  
589 depend on implicit assumptions, which should be articulated.
- 590 • The authors should reflect on the factors that influence the performance of the approach.  
591 For example, a facial recognition algorithm may perform poorly when image resolution  
592 is low or images are taken in low lighting. Or a speech-to-text system might not be  
593 used reliably to provide closed captions for online lectures because it fails to handle  
594 technical jargon.
- 595 • The authors should discuss the computational efficiency of the proposed algorithms  
596 and how they scale with dataset size.
- 597 • If applicable, the authors should discuss possible limitations of their approach to  
598 address problems of privacy and fairness.
- 599 • While the authors might fear that complete honesty about limitations might be used by  
600 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
601 limitations that aren't acknowledged in the paper. The authors should use their best  
602 judgment and recognize that individual actions in favor of transparency play an impor-  
603 tant role in developing norms that preserve the integrity of the community. Reviewers  
604 will be specifically instructed to not penalize honesty concerning limitations.

605 **3. Theory Assumptions and Proofs**

606 Question: For each theoretical result, does the paper provide the full set of assumptions and  
607 a complete (and correct) proof?

608  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647  
648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659  
660  
661  
662

Answer: [Yes]

Justification: Theory is provided in Section 2 and Section 3. Proofs are provide in Appendix A

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. **Experimental Result Reproducibility**

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: All details for reproducing results and experiments can be found through the main text body and appendix. The details include: dataset resource Ramos Zapatero et al. (2023), data processing, model architecture and optimization details, and performance metrics.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
  - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
  - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
  - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701  
702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713

## 5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: The data used in the empirical study is either synthetic or publicly available. The code reproducing all the experiments is attached to the paper.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

## 6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: The paper discusses the experimental setup necessary to understand the results in Section 5. Furthermore, the details of the empirical study are provided in Appendix B.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

## 7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: All the results presented in the paper are averaged over multiple independent runs and the standard deviations are provided along the metrics.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.

- 714 • The factors of variability that the error bars are capturing should be clearly stated (for  
715 example, train/test split, initialization, random drawing of some parameter, or overall  
716 run with given experimental conditions).
- 717 • The method for calculating the error bars should be explained (closed form formula,  
718 call to a library function, bootstrap, etc.)
- 719 • The assumptions made should be given (e.g., Normally distributed errors).
- 720 • It should be clear whether the error bar is the standard deviation or the standard error  
721 of the mean.
- 722 • It is OK to report 1-sigma error bars, but one should state it. The authors should  
723 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis  
724 of Normality of errors is not verified.
- 725 • For asymmetric distributions, the authors should be careful not to show in tables or  
726 figures symmetric error bars that would yield results that are out of range (e.g. negative  
727 error rates).
- 728 • If error bars are reported in tables or plots, The authors should explain in the text how  
729 they were calculated and reference the corresponding figures or tables in the text.

## 730 8. Experiments Compute Resources

731 Question: For each experiment, does the paper provide sufficient information on the com-  
732 puter resources (type of compute workers, memory, time of execution) needed to reproduce  
733 the experiments?

734 Answer: [Yes]

735 Justification: The paper discuss the compute resources and reproducibility in Appendix C.

736 Guidelines:

- 737 • The answer NA means that the paper does not include experiments.
- 738 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,  
739 or cloud provider, including relevant memory and storage.
- 740 • The paper should provide the amount of compute required for each of the individual  
741 experimental runs as well as estimate the total compute.
- 742 • The paper should disclose whether the full research project required more compute  
743 than the experiments reported in the paper (e.g., preliminary or failed experiments that  
744 didn't make it into the paper).

## 745 9. Code Of Ethics

746 Question: Does the research conducted in the paper conform, in every respect, with the  
747 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

748 Answer: [Yes]

749 Justification: The research does conform with the NeurIPS Code of Ethics. The study  
750 presented involves only public or synthetic data, which is freely available online. The  
751 considered models do not impose risks of misuse or dual-use.

752 Guidelines:

- 753 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 754 • If the authors answer No, they should explain the special circumstances that require a  
755 deviation from the Code of Ethics.
- 756 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-  
757 eration due to laws or regulations in their jurisdiction).

## 758 10. Broader Impacts

759 Question: Does the paper discuss both potential positive societal impacts and negative  
760 societal impacts of the work performed?

761 Answer: [Yes]

762 Justification: The paper discusses the broader impact in Appendix E.

763 Guidelines:

- 764 • The answer NA means that there is no societal impact of the work performed.

- 765 • If the authors answer NA or No, they should explain why their work has no societal  
766 impact or why the paper does not address societal impact.
- 767 • Examples of negative societal impacts include potential malicious or unintended uses  
768 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations  
769 (e.g., deployment of technologies that could make decisions that unfairly impact specific  
770 groups), privacy considerations, and security considerations.
- 771 • The conference expects that many papers will be foundational research and not tied  
772 to particular applications, let alone deployments. However, if there is a direct path to  
773 any negative applications, the authors should point it out. For example, it is legitimate  
774 to point out that an improvement in the quality of generative models could be used to  
775 generate deepfakes for disinformation. On the other hand, it is not needed to point out  
776 that a generic algorithm for optimizing neural networks could enable people to train  
777 models that generate Deepfakes faster.
- 778 • The authors should consider possible harms that could arise when the technology is  
779 being used as intended and functioning correctly, harms that could arise when the  
780 technology is being used as intended but gives incorrect results, and harms following  
781 from (intentional or unintentional) misuse of the technology.
- 782 • If there are negative societal impacts, the authors could also discuss possible mitigation  
783 strategies (e.g., gated release of models, providing defenses in addition to attacks,  
784 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from  
785 feedback over time, improving the efficiency and accessibility of ML).

## 786 11. Safeguards

787 Question: Does the paper describe safeguards that have been put in place for responsible  
788 release of data or models that have a high risk for misuse (e.g., pretrained language models,  
789 image generators, or scraped datasets)?

790 Answer: [NA] .

791 Justification: The models considered in the paper do not carry the risks of misuse or dual-use.

792 Guidelines:

- 793 • The answer NA means that the paper poses no such risks.
- 794 • Released models that have a high risk for misuse or dual-use should be released with  
795 necessary safeguards to allow for controlled use of the model, for example by requiring  
796 that users adhere to usage guidelines or restrictions to access the model or implementing  
797 safety filters.
- 798 • Datasets that have been scraped from the Internet could pose safety risks. The authors  
799 should describe how they avoided releasing unsafe images.
- 800 • We recognize that providing effective safeguards is challenging, and many papers do  
801 not require this, but we encourage authors to take this into account and make a best  
802 faith effort.

## 803 12. Licenses for existing assets

804 Question: Are the creators or original owners of assets (e.g., code, data, models), used in  
805 the paper, properly credited and are the license and terms of use explicitly mentioned and  
806 properly respected?

807 Answer: [Yes] .

808 Justification: We cite (Ramos Zapatero et al., 2023) that produced the dataset used in the  
809 study. The dataset is available under the license CC BY 4.0.

810 Guidelines:

- 811 • The answer NA means that the paper does not use existing assets.
- 812 • The authors should cite the original paper that produced the code package or dataset.
- 813 • The authors should state which version of the asset is used and, if possible, include a  
814 URL.
- 815 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 816 • For scraped data from a particular source (e.g., website), the copyright and terms of  
817 service of that source should be provided.

- 818
- 819
- 820
- 821
- 822
- 823
- 824
- 825
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
  - For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
  - If this information is not available online, the authors are encouraged to reach out to the asset's creators.

### 826 13. New Assets

827 Question: Are new assets introduced in the paper well documented and is the documentation  
828 provided alongside the assets?

829 Answer: [NA] .

830 Justification: The paper does not release new assets.

831 Guidelines:

- 832
- 833
- 834
- 835
- 836
- 837
- 838
- 839
- The answer NA means that the paper does not release new assets.
  - Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
  - The paper should discuss whether and how consent was obtained from people whose asset is used.
  - At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

### 840 14. Crowdsourcing and Research with Human Subjects

841 Question: For crowdsourcing experiments and research with human subjects, does the paper  
842 include the full text of instructions given to participants and screenshots, if applicable, as  
843 well as details about compensation (if any)?

844 Answer: [NA] .

845 Justification: The empirical study presented in the paper is conducted on the synthetic or  
846 publicly available data.

847 Guidelines:

- 848
- 849
- 850
- 851
- 852
- 853
- 854
- 855
- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
  - Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
  - According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

### 856 15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human 857 Subjects

858 Question: Does the paper describe potential risks incurred by study participants, whether  
859 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)  
860 approvals (or an equivalent approval/review based on the requirements of your country or  
861 institution) were obtained?

862 Answer: [NA]

863 Justification: The empirical study presented in the paper is conducted on the synthetic or  
864 publicly available data.

865 Guidelines:

- 866
- 867
- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.



868  
869  
870  
871  
872  
873  
874  
875

- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.